

05-25-00

JC711 U.S. PTO
05/23/00

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

Inventorship..... Zizys et al.
Applicant..... Microsoft Corporation
Attorney's Docket No. MS1-517US
Title: Load Simulation Tool for Server Resource Capacity Planning

JC759 U.S. PTO
09/577118
05/23/00

TRANSMITTAL LETTER AND CERTIFICATE OF MAILING

To: Commissioner of Patents and Trademarks,
Washington, D.C. 20231

From: James R. Banowsky (Tel. 509-324-9256; Fax 509-323-8979)
Lee & Hayes, PLLC
421 W. Riverside Avenue, Suite 500
Spokane, WA 99201

The following enumerated items accompany this transmittal letter and are being submitted for the matter identified in the above caption.

1. Specification—title page, plus 32 pages, including 37 claims and Abstract
2. Transmittal letter including Certificate of Express Mailing
3. 4 Sheets Formal Drawings (Figs. 1-4)
4. Return Post Card

Large Entity Status ☒ [x]

Small Entity Status ☐ []

Date: 5/23/00

By:

James R. Banowsky
James R. Banowsky
Reg. No. 37,773

CERTIFICATE OF MAILING

I hereby certify that the items listed above as enclosed are being deposited with the U.S. Postal Service as either first class mail, or Express Mail if the blank for Express Mail No. is completed below, in an envelope addressed to The Commissioner of Patents and Trademarks, Washington, D.C. 20231, on the below-indicated date. Any Express Mail No. has also been marked on the listed items.

Express Mail No. (if applicable) EL624351643

Date: 5/23/00

By:

Lori A. Vierra
Lori A. Vierra

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

**Load Simulation Tool For Server
Resource Capacity Planning**

Inventor(s):
Matt Odhner
Giedrius Zizys

ATTORNEY'S DOCKET NO. MS1-517US

1 **TECHNICAL FIELD**

2 This invention relates to server systems, and more particularly to systems
3 and methods for server resource capacity planning in server systems.
4

5 **BACKGROUND**

6 Capacity planning is forward-looking resource management that allows a
7 computer system administrator to plan for expected changes of system resource
8 utilization and make changes to the system to adequately handle such changes.
9 Server performance and capacity planning is a top concern of computer
10 administrators and business managers. If a lack of proactive and continuous
11 capacity planning procedure leads to unexpected unavailability and performance
12 problems, the downtime that results can be financially devastating to a company
13 that depends heavily on server performance, such as an Internet-based merchant.

14 The importance of superior capacity planning is heightened by the
15 continuous growth in server-dependent companies and potential customers for
16 such companies. Even a solid company that has millions of customers can quickly
17 decline in popularity if it does not increase its resources to handle a constant
18 increase in customers. Excessive downtime of such a company can cause
19 customers to take their business elsewhere.

20 Capacity planning requires both scientific and intuitive knowledge of a
21 server system. It requires in-depth knowledge of the resource being provided and
22 an adequate understanding of future server traffic. The difficulty of the problem
23 has increased by technology in which multiple servers, or server clusters, are
24 employed to handle a network or an Internet website.
25

Current capacity planning methods do not adequately estimate a number of servers having certain resources that a system will need to handle expected loads (requests per second). Therefore, a capacity planning method and system is needed in which a user can provide an expected load that the system needs to handle and receive information on how to increase servers and/or resources to adequately handle that load.

SUMMARY

A method and system for providing capacity planning of server resources is described herein. The methods and systems contemplate using measured data, extrapolation, and a load simulation tool to provide capacity planning results that are more accurate than current schemes. The load simulation tool and its implementation are also described. Server resources for which utilization is calculated are processor utilization, communication bandwidth utilization, memory utilization, and general server utilization.

Utilization is expressed in terms of actual use of the resource in relation to the total amount of resource available for use. For example, processor utilization is expressed as a percentage of procession power utilized for a given load in relation to the total processing power available. Communication bandwidth utilization is expressed as a percentage of an average server throughput per bytes per second in relation to the total communication bandwidth available. Memory utilization is expressed as a percentage of memory required per request times the length of a request queue in relation to the total memory available. General server utilization is expressed as a ratio between a current service rate (number of requests per second served) and the maximum possible service rate (maximum

number of requests the server is capable of serving). This is less specific than showing the processor, bandwidth, and memory utilization, but it is useful for viewing resource constraints that do not fall under the other three categories.

The calculations that are used to derive utilization percentages of server resources require that the maximum load that can be handled by the server cluster (maximum requests / second) be determined. Other methods to estimate this maximum load are described in a related patent application entitled, "Capacity Planning For Server Resources," by Odhner and Zizys, U.S. Patent Application No. 08/_____, filed _____. It is noted that the inventors of the referenced patent application are the same of those of the present application, and that Microsoft Corp. is the assignee of both inventions.

The implementation described herein derives the maximum load of a server cluster by collecting actual server parameter values during operation of the server system. This is accomplished through the use of a filter, such as an Internet Server Application Program Interface (ISAPI) filter, that collects actual server traffic information as data is transmitted to and from the server cluster. In addition, a monitor on each server in the server cluster collects other server parameter values that are used in subsequent calculations.

After the filter and the monitors have collected the required data, a system user selects a client computer from which to run a load simulation tool. The load simulation tool, in effect, replays the data that has been collected from the server cluster, such as the actual requests made to the server, the time intervals at which requests were made, etc. The load simulation tool is then used to increase the load on the system until a maximum service rate that the system can support is found.

There are several ways to calibrate the server load to find the maximum service rate. The number of users from the actual recorded data can be multiplied to simulate a greater number of users, which will increase the load on the system. Another way is to decrease the amount of time between requests, as recorded by the system, which will increase the load on the system. As the load increases, a service rate is monitored. When a further increase in the load does not increase the service rate, the load on the system at that point is considered to be the maximum service rate that can be delivered by the server.

It is noted that the user can create a script manually, instead of replaying the recorded data to calibrate the maximum load, but this will not provide a similarly accurate outcome, since the user in that situation, is required to estimate certain server usage parameters.

After the system is calibrated to find the maximum load that can be handled by the system, the maximum load value is used in subsequent calculations to determine server resource utilization estimates for any number of hypothetical situations. For instance, a user can enter information regarding a particular load that the user wants the current system to handle. The described implementation provides that user with estimates as to the utilization that the specified load will cause for the processor, the memory, the communications bandwidth, and the server in general. Also, the user may want to see how adding or removing a server from a current system will affect the utilization of these server resources. This situation can be adequately determined using the implementation described herein.

Finally, after the user runs the load simulation tool to calibrate the system as to the maximum load and make determinations regarding utilization of server resources, the system provides a plan that recommends any changes in

configuration, if any, that should be made to the system to optimize system performance. These recommendations are stored for each test result, thereby enabling the user to run several tests, and contrast and compare results and recommendations for different situations that the user may expect in the future. The user is thus enabled to adequately plan for future situations.

BRIEF DESCRIPTION OF THE DRAWINGS

A more complete understanding of the various methods and arrangements of the present invention may be had by reference to the following detailed description when taken in conjunction with the accompanying drawings, wherein:

Fig. 1 is an illustration of a prior art server-client system having a server cluster that supports a website on the Internet.

Fig. 2 is a high-level block diagram of a server cluster having a stress simulation tool for capacity planning.

Fig. 3 is a screen shot of a capacity planning worksheet utilized in a capacity planning process using a stress simulation tool.

Fig. 4 is a graph of load vs. processor utilization for a calibrated method of capacity planning.

DETAILED DESCRIPTION

Fig. 1 shows a typical Internet-based server-client system 100. The system 100 includes several clients 104a, 104b, 104c, 104d connected to the Internet 102. A website 106 runs on a server cluster 108 comprised of three servers 110a, 110b, 110c. Although the server-client system 100 is shown operating within an Internet website context, it is noted that the server-client system may operate in any server-

client network context, such as a local area network (LAN) or a wide area network (WAN).

Fig. 2 depicts a server cluster 200 in accordance with the described implementations. The server cluster 200 comprises a primary server 202 having a processor 204 and a monitor 205, a first secondary server 206 having a processor 208 and a monitor 209, and a second secondary processor 210 having a processor 212 and a monitor 213. The monitors are software devices that collect server parameter values while the server cluster 200 is in operation. The server cluster 200 communicates with a master client 214 via a communications connection 216. It is noted that several clients (not shown) may be connected to the server cluster 200. However, only one client is selected by the user to be the master client 214. The master client 214 includes a simulation test program 217. The function of the master client 214 and the simulation test program 217 will be discussed in greater detail below.

The primary server 202 also includes a memory 218 and runs an operating system 220. The operating system 220 provides resource management for primary server 202 resources. The memory 218 of the primary server 202 includes a cluster controller 222, which controls communications between the primary server 202 and the secondary servers 206, 210 and between the server cluster 200 and the network 214. To accomplish this, the cluster controller 222 is provided with a communications program 224.

A capacity planner 226 is included in the cluster controller 222. The function of the capacity planner 226 and its components will be described in greater detail below. Generally, the capacity planner 226 comprises benchmark data 228 in which data collected from the server cluster 200 is stored, a calculation

1 module 230 which stores the equations necessary to derive server resource
2 utilization estimates, and plans 232 which stores recommendations that may be
3 made to improve operational configuration of the server cluster. This file of
4 recommendations is pre-defined by the manufacturer to list all the possible
5 recommendations developed for the server cluster 200. In addition, plans 232 may
6 be updated via a version upgrade or through a connection to the Internet.

7 In addition, the capacity planner 226 includes a user interface 234 and an
8 ISAPI filter 236. The user interface 234 provides areas wherein a user of the
9 server cluster 200 in general and, more specifically, the capacity planner 222 can
10 enter server parameter values and/or a specified load for which the user wants to
11 see server resource utilization and recommendations. The ISAPI filter 236 is used
12 to collect actual server parameter values from the server cluster 200 while the
13 server cluster 200 is operating. It is noted that the filter need not be an ISAPI
14 filter, but can be any type of filter capable of performing the functions listed
15 herein.

16 The capacity planner 222 includes a load simulation tool 238 which is used
17 to construct simulation scripts - such as the simulation test program 217 - that,
18 when run on the master client 214, simulates, plays or replays a server load
19 scenario using actual operating conditions recorded from the server cluster 200.
20 The use of the load simulation tool 238 is described in further detail below.

21 The implementations and functions of the components of the server cluster
22 200 outlined above will become more clear as the discussion progresses with
23 continuing reference to the components of Fig. 2.

24 The server resources that are discussed herein are: (1) processor utilization
25 (also referred to as CPU utilization), wherein the processor utilization for a given

load is expressed as a percentage of total processing power available; (2) memory utilization, expressed as a percentage of total memory available is determined by multiplying the memory required for each request by the number of requests; (3) communication bandwidth utilization, expressed as a percentage of the average throughput per bytes per second in relation to the total communication bandwidth available; and (4) general server utilization, expressed as a ratio between a current service rate (number of requests per second served) and the maximum possible service rate (maximum number of requests the server is capable of serving). The general server utilization is less specific than showing the processor, bandwidth, and memory utilization, but it is useful for viewing resource constraints that do not fall under the other categories.

Fig. 3 shows a screen shot of a user interface 300 for a capacity planning worksheet, wherein the user enters the specified load, for which the user desires to observe the effects on the system of handling such a load. The user is required to manually enter several server parameter values. These server parameter values include: number of servers in the server cluster, available communications bandwidth, server name on which a simulation will be run, client name of the client that will serve as the master test client and execute a simulation script, and the name of the script that will be used to run the simulation.

To begin, the user notifies the server cluster 200 to begin collecting data. The monitors 205, 209, 213 collect data from each server 202, 206, 210. The ISAPI filter 236 collects data for other server parameters, namely for communications-related parameters such as number of incoming requests and average response time.

1 The server resource utilization calculations require knowledge of the
2 maximum load that the server cluster 200 can, theoretically, handle. The
3 implementation described herein is more accurate in deriving the maximum load
4 than any other method described to date.

5 To find this maximum load, actual operating parameters are collected from
6 the server cluster 200 through the monitors 205, 209, 213 and the ISAPI filter 236.
7 The data collected is utilized by the load simulation tool 238 to derive a simulation
8 script that enables the simulation test program 217 on the master client 214 to
9 recreate the server resource utilizations that occurred during the operational
10 period.

11 The simulation is run on only one server, selected by a user via the user
12 interface 300. It is assumed that the primary server 202, and the secondary servers
13 206, 210 are identical. Once the simulation data is derived on one server, the final
14 figures are extrapolated for the total amount of servers in the server cluster. This
15 provides the user with the server resource utilization figures.

16 Although not particularly discussed herein, it is noted that if the servers are
17 not identical, the simulation script can be run on each individual server and then
18 the individual results can be summed to provide the final totals. For discussion
19 purposes, it is assumed that servers 202, 206, 210 are identical.

20 Once a script has been obtained, the user is provided with means to increase
21 the test load on the server to run the script. All the other parameters are the same,
22 so increasing the load will, necessarily, increase the utilization of the server
23 resources.

24 Fig. 4 shows a graph of a load vs. utilization curve 500. For this example,
25 processor utilization is used, though it will be apparent that a similar graph could

1 be constructed for any of the server resource utilization estimates. As the load
2 increases to point 502 on the load axis, the utilization curve 500 reaches a point
3 504 which can be considered to be the maximum load that can be handled by the
4 server 202.

5 The user is may increase the load via the user interface 300, and re-run the
6 script using the higher load value. A situation will arise in which an increase in
7 the load will not result in an increase of the rate at which the load is handled. This
8 is the maximum load 502 which the server can handle. The load (L) at this point
9 is used in the resource utilization estimate calculations below.

10 General server utilization is derived by solving:

$$11 \quad U = \frac{L}{X}$$

12
13
14 wherein:

15 U = general server utilization;

16 L = specified load; and

17 X = maximum load that can be handled by the server cluster 200.
18
19
20
21
22
23
24
25

Processor utilization is derived by solving:

$$U_{CPU} = \frac{a}{e^{b \cdot L}}$$

wherein:

U_{CPU} is processor utilization;

L is the specified load; and

a and b are processor regression constants derived from applying linear regression methodology to several load/utilization (x,y) pairs measured during the test.

Communications bandwidth utilization is derived by solving:

$$U_B = \frac{F_{TCP}}{B} \cdot (c + d \cdot L)$$

wherein:

U_B is communication bandwidth utilization;

F_{TCP} is a transmission overhead factor that, when applied to a certain size page, results in the actual bandwidth necessary to transmit the page;

L is the specified load;

B is the total communication bandwidth available; and

c and d are bandwidth regression constants derived from applying linear regression methodology to several load/utilization (x,y) pairs measured during the test.

The memory utilization is derived by first solving the following equation to determine the number of concurrent connections:

$$N = \frac{L}{(X - L)} + S1 \cdot L$$

wherein:

N is the number of concurrent connections;

L is the specified load;

X is the maximum load that can be handled by the server cluster 200; and

$S1$ is a connection memory factor that is the adjusted average of the incoming connections at different speeds. For example, suppose that the ISAPI filter 236 has measured the following percentages for connection types:

56K: 50%

ADSL: 20% ***question: what relation to screen shot? ISDN? ***

T1: 20%

T3: 10%.

Then $S1$ is the adjusted average of these connection speeds:

56K: $0.5 * 5.6 = 2.8$ KBytes/sec +

ADSL: $0.2 * 30 = 6$ KBytes/sec +

T1: $0.2 * 150 = 30$ KBytes/sec +

T3 $0.1 * 4500 = 450$ KBytes/sec = 488.8 KBytes/sec.

Then $S1 = 488.8$ KBytes/second.

1 The memory utilization is thus derived by solving:

2

$$3 \quad U_M = \frac{N \cdot (M_{TCP} + M_{IISStruct}) + M_{OS} + M_{IIS}}{4 \quad M}$$

5 wherein:

6 U_M is memory utilization;

7 N is the number of concurrent connections;

8 M_{TCP} is an amount of memory for TCP buffers (32 KB per connection);

9 M_{IIS} is the amount of memory required by a server communication program
10 (50 MB for IIS);

11 $M_{IISStruct}$ is the amount of memory necessary to support communications
12 program data structures associated with each connection (50 KB per connection
13 for IIS);

14 M_{OS} is the amount of memory required by a server operating system (64
15 MB for Windows® NT by Microsoft® Corp.) and

16 M is the amount of total memory available.

17 It is noted that some figures have been used that are specific to IIS, the
18 communications program 224 used for purposes of this discussion. However, it is
19 noted that these numbers may be different for different communications programs.

20

21 Conclusion

22 The described implementations advantageously provide for capacity
23 planning for a server-client system and, particularly, to a server cluster within a
24 server-client system. The load simulation tool is an extremely accurate tool for
25 determining the maximum load handled by a server. The maximum load can then

1 be substituted into the server resource estimate equations to give accurate server
2 resource utilization results.

3 Although the invention has been described in language specific to structural
4 features and/or methodological steps, it is to be understood that the invention
5 defined in the appended claims is not necessarily limited to the specific features or
6 steps described. Rather, the specific features and steps are disclosed as preferred
7 forms of implementing the claimed invention.

1 **CLAIMS**

2 1. A method for deriving server resource utilization estimates for a
3 server cluster, the method comprising:

4 recording server cluster data during operation of the server cluster, at least
5 some of the server cluster data indicating server resource parameter values;

6 using a load simulation tool that, using the recorded data, determines a
7 maximum load that can be handled by the server cluster;

8 specifying a load to be handled by the server cluster; and

9 deriving server resource utilization estimates corresponding to the specified
10 load.

11
12 2. The method as recited in claim 1, further comprising:

13 displaying the server resource utilization estimates; and

14 recommending a plan to optimize processing of the specified load.

15
16 3. The method as recited in claim 2, wherein the plan recommends a
17 change in the hardware configuration of the server cluster.

18
19 4. The method as recited in claim 1, wherein the maximum load, the
20 recorded values, the specified load, and the server resource utilization estimates
21 are stored in non-volatile memory.
22
23
24
25

1
2 **5.** The method as recited in claim 1, wherein the using a load simulation
3 tool comprises:

4 creating a test script from the recorded values;

5 running the test script on a master client to simulate load and server
6 resource utilization conditions that existed on a server when the recorded values
7 were recorded; and

8 increasing the load on the server, when the test script is running, until a
9 maximum load that can be handled by the server is obtained.
10

11 **6.** The method as recited in claim 5, wherein the increasing the load on
12 the server further comprises multiplying the number of users utilizing the server
13 cluster when the recorded values were recorded, thereby multiplying the resources
14 utilized by the users.
15

16 **7.** The method as recited in claim 5, wherein the increasing the load on
17 the server further comprises decreasing the amount of time between user requests,
18 thereby increasing the resources utilized by the users.
19
20
21
22
23
24
25

1 **8.** The method as recited in claim 5, wherein the increasing the load on
2 the server until a maximum load that can be handled by the server is obtained,
3 further comprises:

4 observing a service rate exhibited by the server on which the simulation is
5 being performed; and

6 recognizing that the maximum load has been obtained when an increase in
7 the load does not increase the service rate.

8
9 **9.** The method as recited in claim 5, wherein:

10 the server cluster contains a set of identical servers;

11 running the test script run on the master client simulates server cluster
12 operation on only one of the servers of the server cluster; and

13 the method further comprises extrapolating the results obtained on the one
14 server using the number of servers in the set of identical servers to obtain the
15 maximum load that can be handled by the server cluster.

16
17 **10.** The method as recited in claim 5, wherein:

18 the server cluster contains a set of non-identical servers;

19 running the test script run on the master client further comprises running
20 the test script on each of the non-identical servers in the server cluster; and

21 the method further comprises summing the results obtained from each non-
22 identical server in the cluster to obtain the maximum load that can be handled by
23 the server cluster.

1 16. The method as recited in claim 1, wherein the server resource
2 utilization comprises memory utilization, the method further comprising:

3
4 deriving memory utilization by solving:

5
6
$$U_M = \frac{N \cdot (M_{TCP} + M_{IISStruct}) + M_{OS} + M_{IIS}}{M}$$

7

8
9 wherein N is a total number of concurrent connections derived by solving:

10
$$N = \frac{L}{(X - L)} + S1 \cdot L$$

11

12 wherein: U_M is memory utilization; M_{TCP} is a an amount of memory
13 necessary to support the connections for communications; $M_{IISStruct}$ is the amount
14 of memory necessary to support data structures associated with each connection;
15 M_{OS} is the amount of memory required by a server operating system; M_{IIS} is the
16 amount of memory required by a server communication program; M is the total
17 amount of memory available; L is the specified load; X is the maximum load that
18 can be handled by the server cluster; and $S1$ is a connection memory factor that is
19 the adjusted average of the incoming connections at different speeds.
20
21
22
23
24
25

1
2 **21.** The method recited in claim 1, wherein the server resource
3 utilization is processor utilization, the method further comprising:

4 finding a functional dependency approximation between processor
5 utilization and load;

6 transforming functional dependency into linear form by using logarithmic
7 transformation;

8 deriving first and second processor regression constants using linear
9 regression methodology;

10 dividing the first processor regression constant by e to the power of the
11 product of the second processor regression constant and the specified load to
12 obtain the processor utilization estimate.

13
14 **22.** The method as recited in claim 1, wherein the server resource
15 utilization is communication bandwidth utilization, the method further comprising:

16 finding a functional dependency approximation between communication
17 bandwidth utilization;

18 transforming functional dependency into linear form by using logarithmic
19 transformation;

20 deriving first and second bandwidth regression constants using linear
21 regression methodology;

22 deriving a transmission overhead factor that, when applied to a certain size
23 web page, results in the actual capacity necessary to transmit the web page;

24 deriving a weighted communication overhead factor by dividing the
25 transmission overhead factor by the available communication bandwidth;

1 deriving an adjusted communication load by adding the first bandwidth
2 regression constant to the product of the specified load and the second bandwidth
3 regression constant; and

4 determining the communication bandwidth utilization estimate by
5 multiplying the weighted communication overhead factor by the adjusted
6 communication load.

7
8 **23.** The method as recited in claim 1, wherein the server resource
9 utilization is memory utilization, the method further comprising:

10 deriving a connection memory factor that is the adjusted average of the
11 incoming connections at different speeds;

12 deriving a weighted connection memory factor by multiplying the
13 connection memory factor by the specified load;

14 deriving a page load ratio by dividing the specified load by the difference of
15 the maximum load value and the specified load;

16 deriving a total number of concurrent connections by adding the weighted
17 connection memory factor and the page load ratio; and

18 deriving a gross memory utilization by multiplying the total number of
19 concurrent connections by the sum of the amount of memory necessary to support
20 each connection for communications and the amount of memory necessary to
21 support data structures associated with each connection, and adding the amount of
22 memory required by a server operating system and the amount of memory
23 required by the server communication program; and

24 deriving the memory utilization estimate by dividing the gross memory
25 utilization by total memory available.

1
2 **24.** The method as recited in claim 1, wherein the server resource
3 utilization is general server utilization, the method further comprising:

4 dividing the specified load by the maximum load to derive the general
5 server utilization estimate.
6

7 **25.** One or more computer-readable media having computer-readable
8 instructions thereon which, when executed by one or more computers, cause the
9 computers to implement the method of claim 1.
10

11 **26.** A simulation tool for use in determining server resource utilization
12 estimates in a server cluster having one or more servers, the load simulation tool
13 comprising:

14 a user interface configured to receive data input from a user;

15 at least one filter or monitor configured to record operational data from one
16 or more of the servers in the server cluster;

17 the simulation tool being configured to create a test script from the recorded
18 data and the received data, and to run the test script from a master client connected
19 to the server cluster to simulate load and other server conditions that existed when
20 the operational data was recorded; and

21 the user interface being further configured to display utilization of server
22 resources during the running of the test script.
23
24
25

1 27. The simulation tool as recited in claim 26, wherein the simulation
2 tool being configured to create the test script is further configured to allow the user
3 to create the test script to increase the load on the server on which the simulation is
4 running and observe the effect of such an increase in load on the server resource
5 utilization displays.

6
7 28. A system, comprising:
8 a server cluster having one or more servers, one of which is a primary
9 server that controls the operation of the server cluster;
10 a cluster controller resident in memory on the primary server of the server
11 cluster, the cluster controller controlling communications between the primary
12 server and secondary servers, if any, and between clients and the server cluster;
13 an operating system resident in the memory of the primary server;
14 a communications program within the cluster controller to provide
15 communications capability for the system;
16 a filter to collect server data indicating certain operating parameters for the
17 server cluster;
18 a monitor on each server in the server cluster to collect server data
19 indicating certain operating parameters for the server cluster;
20 a user interface to collect data input by a user;
21 a capacity planner within the cluster controller configured to utilize the
22 collected data to derive one or more server resource utilization estimates for server
23 resources to determine how handling a specified load will affect the utilization of
24 the server resources, and to produce a plan recommending changes to be made to
25 the server cluster to adequately accommodate the specified load; and

1 a load simulation tool configured to use the collected data to create a
2 simulation script that, when run on a master client, simulates the operation of the
3 server cluster system to allow the user to find the maximum load that the server
4 cluster can handle; and

5 wherein the maximum load obtained through the use of the load simulation
6 tool is utilized in the derivation of the one or more server resource utilization
7 estimates.

8
9 **29.** The system as recited in claim 28, wherein the filter is an ISAPI
10 filter.

11
12 **30.** The system as recited in claim 28, wherein the collected data and the
13 plans are stored in the memory.

14
15 **31.** The system as recited in claim 28, wherein the simulation script is
16 run from a master client connected to the server cluster, and wherein the
17 simulation is performed on only one server of the server cluster.

18
19 **32.** The system as recited in claim 31, wherein the load simulation tool
20 is further configured to extrapolate results from the simulation on one server in the
21 server cluster to derive results for the total number of servers in the server cluster.

1
2 37. The system as recited in claim 28, wherein the server resource
3 utilization derived by the capacity planner comprises communication bandwidth
4 utilization, and the capacity planner is further configured to derive communication
5 bandwidth utilization by solving:

$$U_M = \frac{N \cdot (M_{TCP} + M_{IISStruct}) + M_{OS} + M_{IIS}}{M}$$

6
7
8
9 wherein N is a total number of concurrent connections derived by solving:

$$N = \frac{L}{(X - L)} + S1 \cdot L$$

10
11
12
13
14 wherein: U_M is memory utilization; M_{TCP} is a an amount of memory
15 necessary to support the connections for communications; $M_{IISStruct}$ is the amount
16 of memory necessary to support data structures associated with each connection;
17 M_{OS} is the amount of memory required by a server operating system; M_{IIS} is the
18 amount of memory required by a server communication program; M is the total
19 amount of memory available; L is the specified load; X is the maximum load that
20 can be handled by the server cluster; and $S1$ is a connection memory factor that is
21 the adjusted average of the incoming connections at different speeds.
22
23
24
25

1 **ABSTRACT**

2 A methods and systems for capacity planning of server resources are
3 described wherein a load simulation tool is used to use actual data gathered from a
4 server cluster during operation to simulate server cluster operation in which the
5 load (requests per second) can be increased, and the effects on the utilization of
6 resources can be observed. Plans containing recommendations are then presented
7 to a system user so the user can make decisions necessary regarding whether to
8 change configuration hardware to meet expected load increases in the future.
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Fig. 1
Prior Art

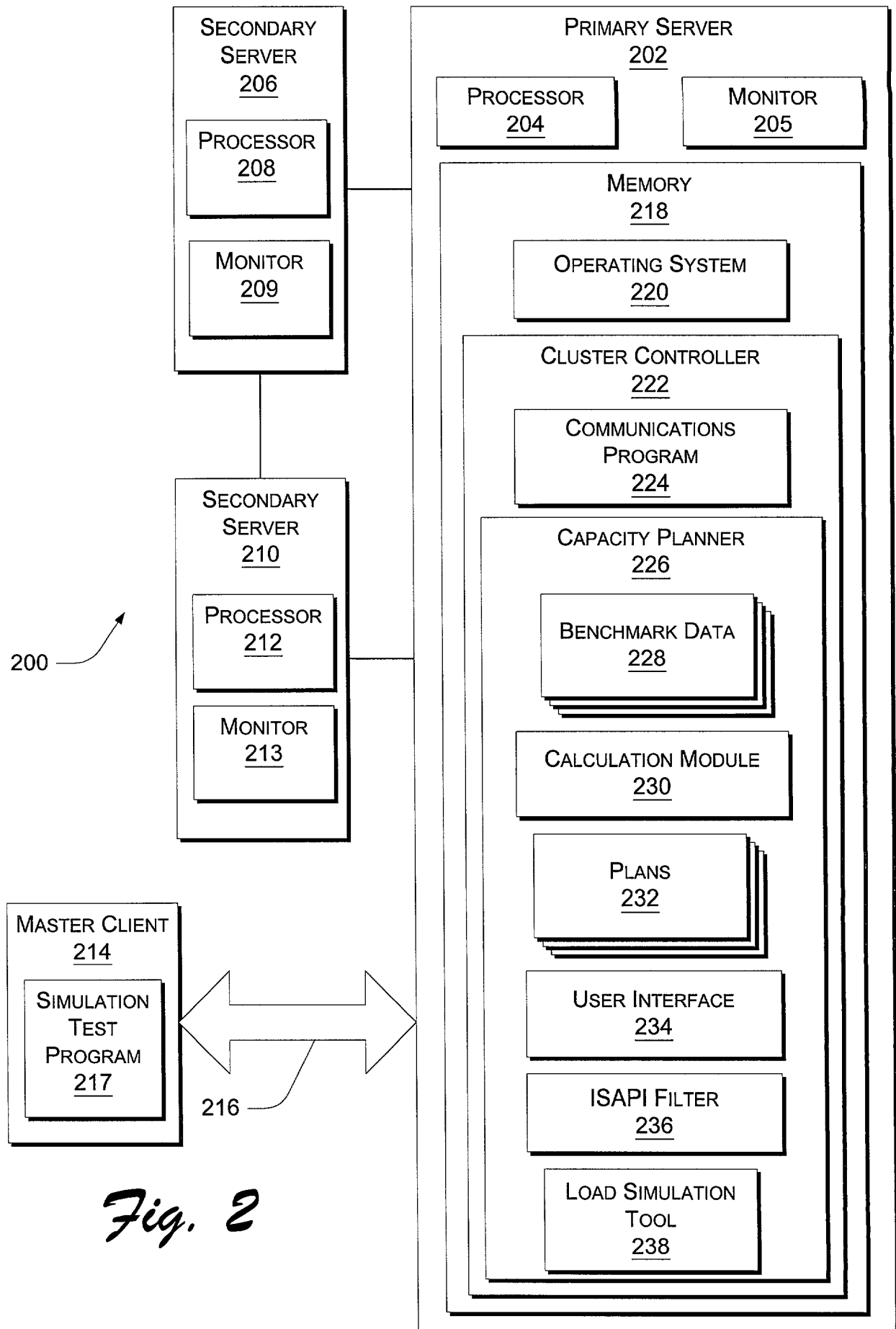


Fig. 2

APPCENTER CAPACITY PLANNING WORKSHEET

CALIBRATION ANALYSIS

MACHINE INFORMATION

CONFIGURE CALIBRATION

CALIBRATION SCRIPT: PG MIX 1

RUN NOW

CONFIGURE CALIBRATION TEST

SERVER NAME: WEB02

MASTER TEST CLIENT NAME: CL01T

TEST SCRIPT: SCRIPT022A

CALCULATE PROJECTED CAPACITY

REQUESTS PER SECOND 37

NUMBER OF SERVERS 8

AVAILABLE BANDWIDTH: T1 (1.54 Mb/Sec)

HELP

WORKSHEET RESULTS

CALIBRATED DATA

GEN.		77%
PROC.		40%
B/W		98%
MEM.		64%

RECOMMENDATIONS

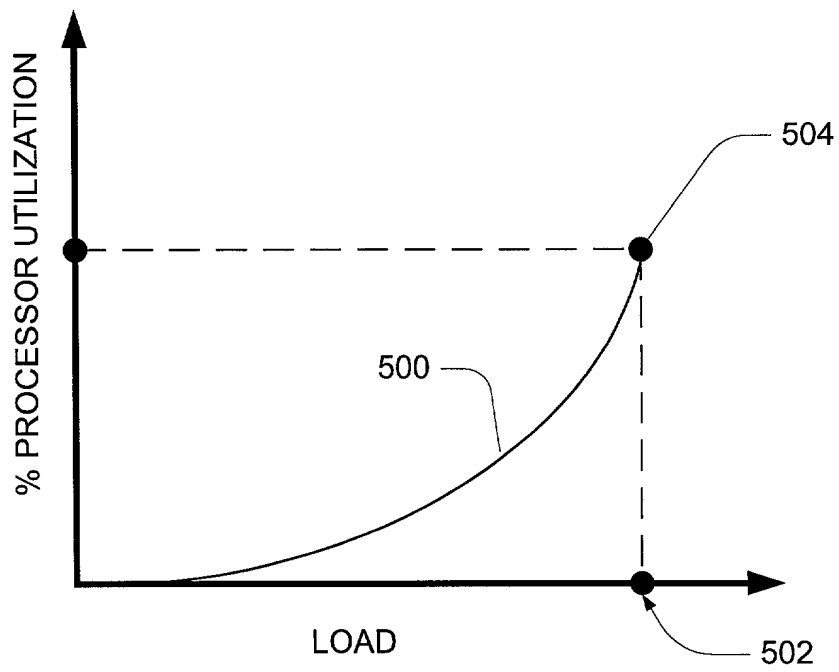
STATUS: CALIBRATION RUNNING....

<<PREVIOUS

NEXT >>

300

Fig. 3



CALIBRATION RESULTS

Fig. 4